

¹ BEYOND THE SCOPE: WILL HUMANOID ROBOTS
² DEFINE THE FUTURE OF ENDOSCOPY?

³ **Yichen Wang¹, Yuting Huang¹, Mario El Hayek², Vivek Kumbhari¹, and**
⁴ **Michael B. Wallace¹**

⁵ ¹Division of Gastroenterology and Hepatology, Mayo Clinic, Jacksonville, FL,
⁶ USA

⁷ ²Department of Medicine, Mayo Clinic, Jacksonville, FL, USA

⁸ **Correspondence**

⁹ Michael B. Wallace, MD, MPH, Mayo Clinic, 4500 San Pablo Rd, Jacksonville, FL 32224

¹⁰ For the past half-century, the technological trajectory of gastrointestinal endoscopy has been fo-
¹¹ cused on visualization. Innovations in optics and high-definition sensors were designed to present
¹² a clearer picture to the human eye, placing the task of interpretation solely on the physician. Re-
¹³ cent years have seen the integration of artificial intelligence (AI), in the form of computer-aided
¹⁴ detection (CADe) and diagnosis (CADx).¹ Systems such as GI Genius have demonstrated the abil-
¹⁵ ity to reduce adenoma miss rates by acting as a "second observer"² However, these advancements
¹⁶ remain fundamentally passive; they reside on the screen, disconnected from the physical activity
¹⁷ of the procedure. The endoscopist still bears the total cognitive and physical load of navigating

18 the tortuous anatomy of the colon, creating a bottleneck where fatigue and variable motor skills
19 contribute to significant inter-observer variability in procedural quality^{3,4}.

20 The discipline is now standing at the precipice of the epoch of embodied AI. This new paradigm
21 represents the integration of advanced cognitive architectures, specifically vision-language-action
22 (VLA) models, with physical robotic actuation. It marks the transition from systems that merely
23 "see" to systems that "act." This shift is not merely incremental; it is a fundamental reimagining
24 of what an endoscopic system can be. Embodied AI envisions a robotic agent capable of per-
25 ceiving the complex, deformable environment of the gut, reasoning about anatomical landmarks
26 in semantic terms, and autonomously executing kinematic maneuvers. This review analyzes the
27 theoretical underpinnings of this convergence, exploring the specialized hardware of soft robotics
28 and the emerging, disruptive potential of humanoid robots as autonomous operators.

29 The brain of an embodied agent requires capabilities far exceeding old convolutional neural net-
30 works. The field is adopting foundational vision-language models, to bridge medical knowledge
31 and vision signals. Architectures like Endo-FM⁵ and Endo-DINO⁶ provide pre-trained founda-
32 tional models for downstream tasks. The most significant leap is the Endoscopic Vision-Language-
33 Action (EndoVLA) model.⁷ EndoVLA is an end-to-end framework where a live video feed and
34 a human text prompt generate direct control commands. This architecture solves the "semantic
35 gap" in robotics by translating high-level intent, such as "resect that polyp," into low-level motor
36 signals. To domesticate general-purpose VLMs for this safety-critical task, researchers employ
37 a Dual-Phase Fine-Tuning strategy.⁷ The model first undergoes Supervised Fine-Tuning to learn
38 the "grammar" of robotic kinematics, followed by Reinforcement Fine-Tuning to hone precision
39 through interaction.

40 A sophisticated cognitive engine is futile without a physical body capable of safe execution. One
41 dominant approach to this challenge is the redesign of the instrument itself. Traditional semi-
42 rigid endoscopes are prone to causing painful looping; consequently, the hardware for autonomous
43 GI endoscopy is evolving toward soft robotics, prioritizing compliance and safety.^{8,9} A paradigm-

⁴⁴ shifting example is the spider-inspired magnetic soft robot developed at the University of Macau¹⁰.
⁴⁵ Modeled after the golden wheel spider, this device utilizes a magneto-active elastomer matrix to
⁴⁶ switch between rolling and climbing locomotion modes. Unlike passive capsules, this robot can
⁴⁷ actively climb over haustral folds and navigate upside down via magnetic anchoring. Simultane-
⁴⁸ ously, Magnetic Capsule Endoscopy is maturing from teleoperated systems to autonomous plat-
⁴⁹ forms. Projects like MAGIC-AIM integrate VLM architectures with magnetic guidance, allowing
⁵⁰ the capsule to identify the pylorus and autonomously execute the traversal maneuver^{7,11}. This au-
⁵¹ tomaton ensures a complete map of the gastric mucosa is generated regardless of operator fatigue.

⁵² While soft robotics reinvent the tool, a parallel lineage proposes replacing the operator. The rise of
⁵³ general-purpose humanoid robots offers a tantalizing alternative: an autonomous agent that walks
⁵⁴ into the endoscopy suite, picks up the existing standard colonoscope, and performs the procedure
⁵⁵ using the same ergonomics designed for human hands.

⁵⁶ Backward compatibility drives this approach. Hospitals worldwide have billions of dollars invested
⁵⁷ in standard flexible endoscopes. Soft robotic solutions require abandoning this infrastructure for
⁵⁸ expensive, proprietary magnetic guidance systems. A humanoid robot, conversely, is a "drop-in"
⁵⁹ replacement. It utilizes the current capital equipment, navigating the colonoscope via the control
⁶⁰ wheels and shaft torque just as a human fellow would. Recent feasibility studies have demon-
⁶¹ strated that humanoids can effectively manipulate laparoscopic tools via teleoperation, suggesting
⁶² that the transition to flexible endoscopy is a matter of software refinement rather than hardware
⁶³ impossibility¹².

⁶⁴ The cold hard truth is we are stalled by clinical and technical inertia, primarily the "Dexterity Gap."
⁶⁵ We have spent decades refining the endoscope as a tool for humans, yet we now face a reality
⁶⁶ where our current hardware is a legacy anchor. The industry is hooked on high-margin, semi-
⁶⁷ disposable scopes that are fundamentally incompatible with the precision of high-torque robotic
⁶⁸ motors⁹ If we continue to ignore the friction between 20th-century ergonomics and 21st-century
⁶⁹ AI, we risk creating a perpetual pilot phase that never reaches the patient's bedside. Hospitals

70 are loath to cannibalize their existing infrastructure, yet they complain about a shortage of hands
71 to hold the scope. The human hand is a marvel of compliance and sensory feedback, capable of
72 sensing the subtle resistance of a loop forming in the sigmoid colon through the shaft of the scope.
73 Current commercial humanoid hands are often rigid, under-actuated grippers designed for gross
74 manipulation, not the delicate interplay of torque and translation required in endoscopy¹³. For
75 a humanoid to be viable, it requires reliable, dexterous hands equipped with high-fidelity tactile
76 sensors at a low price point. The robot must be able to manipulate the angulation wheels with its
77 thumb while simultaneously torquing the shaft with its right hand, a bimanual coordination task
78 that remains an open challenge in robotics control.

79 Furthermore, there is a critical scarcity of training data known as the "proprioceptive void." While
80 we have millions of hours of endoscopic video, we have virtually zero data recording the actions
81 of the endoscopist's hands synchronized with that video^{2,9}. We know what the colon looks like,
82 but we do not know what the doctor's hands were doing to produce that view. To train a humanoid
83 EndoVLA, we must bridge this gap by instrumenting expert endoscopists with motion-capture
84 gloves or sensors during live procedures, creating massive "Vision-Action" datasets. Once this
85 data exists, the potential for scaling is boundless. A model trained on the collective motor skills of
86 the world's best endoscopists could be deployed to a humanoid in a rural clinic, instantly granting
87 it the dexterity of a master clinician.

88 The humanoid approach warrants optimism. General robotics outpaces specialized medical tools
89 due to massive tech investment. As costs fall, humanoids will address the global gastroenterologist
90 shortage. A future where robots perform screening colonoscopies 24/7 is a converging probability.
91 Economic incentives are powerful. As reimbursement fluctuates, a robotic workforce becomes
92 an existential necessity. Continuous operation eliminates the 'off-hours' penalty where outcomes
93 suffer due to staffing. A robot does not sleep, and its hand never shakes at the end of a shift. This
94 consistency is the true promise of the revolution.

95 The integration of the brain and body occurs in the navigation loop. FM and VLM architectures

96 provide the reasoning required for autonomous tracking through language-conditioned tasks. By
97 processing a prompt and video feed to output discrete motion commands, these systems offer Zero-
98 Shot Generalization, allowing the tracking of novel objects like surgical clips without retraining⁷.

99 Generative AI necessitates rigorous safety. Researchers are employing Safe Reinforcement Learn-
100 ing using Control Barrier Functions that act as mathematical virtual walls, analyzing actions
101 against safety constraints to prevent perforation.¹⁴ To bridge the Sim-to-Real gap, agents are
102 trained in high-fidelity simulators using domain randomization, ensuring robust performance de-
103 spite the chaotic textures of a living patient.

104 The convergence of embodied AI and GI endoscopy signifies a fundamental restructuring of gas-
105 troenterological care. The technological pillars are in place: Vision-Language Models have solved
106 the semantic gap, enabling robots to understand medical intent; and two distinct physical lineages
107 are vying for dominance. On one hand, soft robotics offers a safety-first reinvention of the instru-
108 ment; on the other, humanoid robotics offers a scalable, backward-compatible reinvention of the
109 operator. As we look toward the next decade, the standard of care is poised to shift from manual
110 operation to supervisory control. Future screening protocols may involve patients being treated
111 by a humanoid agent that wields standard instruments with the precision of a machine and the
112 adaptability of a human. This evolution promises to democratize high-quality care, reducing the
113 physical burden on physicians and ensuring that every patient receives a procedure defined by the
114 collective intelligence of the field rather than the variability of a single operator. The era of the
115 passive tool is ending; the era of the embodied partner has begun.

116 For practicing endoscopists, a "robot-ready room" can log wheel rotations, torque, and finger kine-
117 matics during routine lists to generate vision–action corpora without slowing workflow. Health
118 systems absorbing the screening-age drop to 45 years could designate supervised data-collection
119 blocks and trial overnight humanoid-run screening lines with remote oversight to clear backlogs¹⁵.
120 These experiments demand audit trails, explainable policies, and governance that preserve human
121 accountability even as robots offload repetitive motor work¹⁶.

¹²² **Keywords:** Artificial Intelligence, Humanoid Robots, Soft Robotics, Vision-Language Models

¹²³ ABBREVIATIONS

¹²⁴ AI, Artificial Intelligence; CADe, Computer-Aided Detection; CADx, Computer-Aided Diagnosis;
¹²⁵ EndoVLA, Endoscopic Vision-Language-Action; FM, Foundation Model; GI, Gastrointestinal;
¹²⁶ VLM, Vision-Language Model.

¹²⁷ References

¹²⁸ References

- ¹²⁹ 1. Chan HP, Samala RK, Hadjiiski LM, and Sahiner C. Computer-aided diagnosis in the era of deep learning. *Medical Physics* 2020;47:e218–e227.
- ¹³⁰ 2. Wallace MB, Sharma P, Bhandari P, et al. Impact of Artificial Intelligence on Miss Rate of Colorectal Neoplasia. *Gastroenterology* 2022;163:295–304.
- ¹³¹ 3. Chan MY, Cohen H, and Spiegel BM. Fewer polyps detected by colonoscopy as the day progresses at a Veteran’s Administration teaching hospital. *Clinical Gastroenterology and Hepatology* 2009;7:1217–23.
- ¹³² 4. Kaminski MF, Regula J, Kraszewska E, et al. Quality indicators for colonoscopy and the risk of interval cancer. *New England Journal of Medicine* 2010;362:1795–803.
- ¹³³ 5. Wang Y, Zhang Z, Liu H, et al. Endo-FM: Foundation Model for Endoscopy Video Analysis via Large-scale Pre-training. *arXiv preprint arXiv:2306.16741*. 2023. URL: <https://arxiv.org/abs/2306.16741>.

- 141 6. Liu Y, Wang Y, Zhang Z, Chen H, and Wu J. Endo-DINO: Self-Supervised Pre-training for
142 Endoscopy with Vision Transformers. arXiv preprint arXiv:2501.05488. 2025. URL: <https://arxiv.org/abs/2501.05488>.
- 143
- 144 7. Wang R, Hong Y, He Y, Wu J, and Yuan W. EndoVLA: Progressive Endoscopic Vision-
145 Language-Action Model. In: *Medical Image Computing and Computer Assisted Intervention – MICCAI 2024*. Springer, 2024:446–56. DOI: 10.1007/978-3-031-72083-3_42. URL:
146 https://doi.org/10.1007/978-3-031-72083-3_42.
- 147
- 148 8. Shah SG, Brooker JC, Thapar C, Williams CB, and Saunders BP. Patient pain during colonoscopy:
149 an analysis using real-time magnetic endoscope imaging. *Endoscopy* 2002;34:435–40.
- 150
- 151 9. Martin JW, Scaglioni B, Norton JC, et al. Enabling the future of colonoscopy with intelligent
and autonomous magnetic manipulation. *Nature Machine Intelligence* 2020;2:595–602.
- 152
- 153 10. Wang Z, Shi K, and Xu Q. A Millimeter-Scale Magnetic Soft Robot With Spider-Inspired
Multi-Modal Locomotion and Object Manipulation. *IEEE/ASME Transactions on Mecha-
tronics* 2022;27:4676–87.
- 154
- 155 11. Hong Y, Zhang Y, and Meng MQH. Autonomous pylorus traversal in magnetic capsule en-
156 doscopy via deep reinforcement learning. *IEEE Transactions on Robotics* 2021;37:1094–
157 108.
- 158
- 159 12. Wang R, Hong Y, He Y, Wu J, and Yuan W. A Feasibility Study on Humanoid Robot Tele-
operation for Laparoscopic Surgery. arXiv preprint. 2024.
- 160
- 161 13. Piazza C, Grioli G, Catalano MG, and Bicchi A. A Century of Robotic Hands. *Annual Review
of Control, Robotics, and Autonomous Systems* 2019;2:1–32.
- 162
- 163 14. Hong Y, Wang R, He Y, Wu J, and Yuan W. Safe Reinforcement Learning for Autonomous
Flexible Endoscopy. In: *2023 IEEE International Conference on Robotics and Automation (ICRA)*. 2023:9460–6. DOI: 10.1109/ICRA48891.2023.10161476. URL: <https://doi.org/10.1109/ICRA48891.2023.10161476>.
- 164
- 165

- ¹⁶⁶ 15. Crockett SD and Ladabaum U. Potential Effects of Lowering Colorectal Cancer Screening
¹⁶⁷ Age to 45 Years on Colonoscopy Demand, Case Mix, and Adenoma Detection Rate. Gas-
¹⁶⁸ troenterology 2022;162:984–986.e5.
- ¹⁶⁹ 16. Topol EJ. High-performance medicine: the convergence of human and artificial intelligence.
¹⁷⁰ Nature Medicine 2019;25:44–56.